

Importance Weighting for Aligning Language Models under Deployment Distribution Shift

Thanawat Lodkaew¹ Tongtong Fang² Takashi Ishida^{3,1} Masashi Sugiyama^{3,1}

¹The University of Tokyo ²The Institute of Statistical Mathematics ³RIKEN AIP



I. Summary

- Motivation.** Training and deployment objectives often differ. For example, models are trained for **helpfulness** but deployed for **harmlessness**, creating a *deployment distribution shift*.
- Key assumption.** Within the training dataset, some instances are *useful* (**relevant**), such as those containing helpful and harmless responses, for optimizing performance under the deployment distribution. In contrast, others are *not useful* (**irrelevant**), such as those that are helpful but harmful responses.
- Method.** Inspired by [1], we propose an *importance weighting* (IW) method tailored for *direct preference optimization* (DPO) [2], IW-DPO, to mitigate this distribution shift by estimating *importance weights* through density ratio estimation between training and validation data, **upweighting relevant** instances and **downweighting irrelevant** ones to better align with the deployment distribution.
- Results.** Experimental results under various distribution shift scenarios using multiple datasets demonstrate the effectiveness of our approach, with approximately 4% overall win rate improvement over the standard DPO.

II. Deployment Distribution Shift

i. Definition

The **deployment environment** (**deployment dist.**) changes in ways not reflected in the **training dataset** (**training dist.**) due to *changes in end-user behavior, preferences, etc.*

$$p_{\text{tr}}(x, y_1, y_2, b) \neq p_{\text{te}}(x, y_1, y_2, b)$$

ii. Factors of distribution shift

$$p(x, y_1, y_2, b) = p(x)p(y_1, y_2 | x)p(b | x, y_1, y_2)$$

①

②

③

$$p_{\text{tr}}(x) \neq p_{\text{te}}(x)$$

Prompt

$$p_{\text{tr}}(y_1, y_2 | x) \neq p_{\text{te}}(y_1, y_2 | x)$$

Response

$$p_{\text{tr}}(b | x, y_1, y_2) \neq p_{\text{te}}(b | x, y_1, y_2)$$

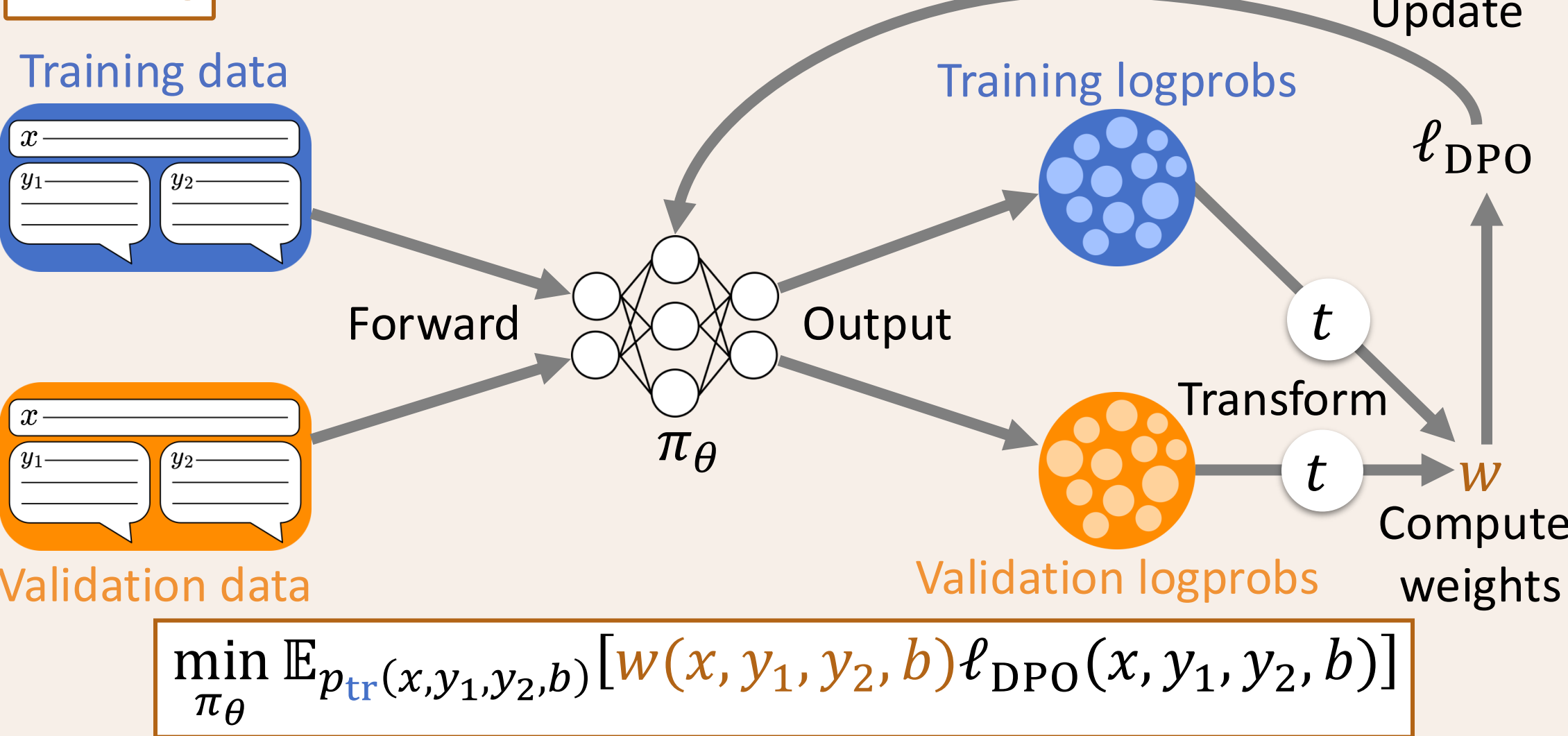
Preference label

iii. Distribution shift types

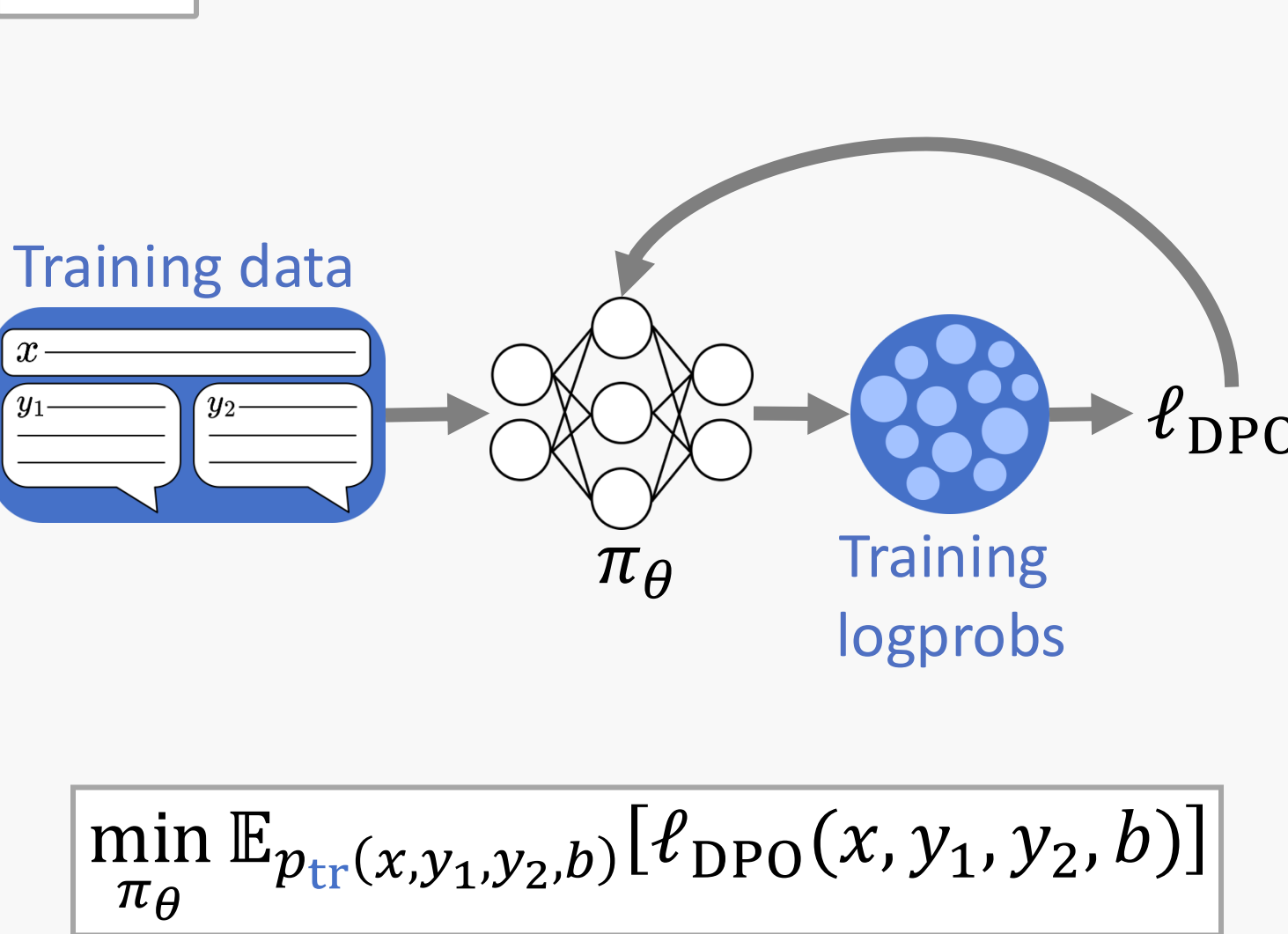
Type of shift	Factor		
	①	②	③
a No shift			
b Full shift	✓	✓	✓
c Prompt shift	✓		
d Response shift		✓	
e Preference label shift			✓
f Prompt + response shift	✓	✓	
g Prompt + preference label shift	✓		✓
h Response + preference label shift		✓	✓

III. Importance Weighted Direct Preference Optimization

IW-DPO

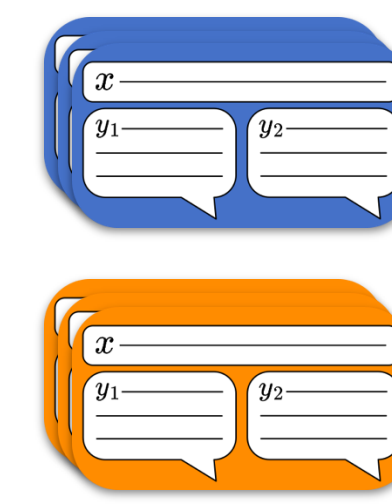


DPO



ii. Problem setting

Training and validation datasets are available, with the constraint that $N_v \ll N_{\text{tr}}$



$$D_{\text{tr}} = \{(x^{\text{tr},i}, y_1^{\text{tr},i}, y_2^{\text{tr},i}, b^{\text{tr},i})\}_{i=1}^{N_{\text{tr}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(x, y_1, y_2, b)$$

$$D_v = \{(x^{\text{v},i}, y_1^{\text{v},i}, y_2^{\text{v},i}, b^{\text{v},i})\}_{i=1}^{N_v} \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(x, y_1, y_2, b)$$

Goal is to *optimize for the test distribution*

$$J(\pi_\theta) = \mathbb{E}_{p_{\text{te}}(x, y_1, y_2, b)}[\ell_{\text{DPO}}(x, y_1, y_2, b)]$$

iii. Definition of importance weight and training objective

Assume support of the training distribution covers that of the test distribution $\text{supp}(p_{\text{te}}) \subseteq \text{supp}(p_{\text{tr}})$

$$w^*(x, y_1, y_2, b) = p_{\text{te}}(x, y_1, y_2, b) / p_{\text{tr}}(x, y_1, y_2, b)$$

$$J(\pi_\theta) = \mathbb{E}_{p_{\text{tr}}(x, y_1, y_2, b)}[w^*(x, y_1, y_2, b) \ell_{\text{DPO}}(x, y_1, y_2, b)] = J_{\text{tr}}(\pi_\theta, w^*)$$

Empirical training objective

$$\hat{J}(\pi_\theta) = \frac{1}{N_{\text{tr}}} \sum_{i=1}^{N_{\text{tr}}} w^{\text{tr},i} \ell_{\text{DPO}}(x^{\text{tr},i}, y_1^{\text{tr},i}, y_2^{\text{tr},i}, b^{\text{tr},i})$$

iv. Importance weight estimation

Transformation function $t: (x, y_1, y_2, b) \mapsto z$

Transformed data

$$Z_{\text{tr}} = \{t(x^{\text{tr},i}, y_1^{\text{tr},i}, y_2^{\text{tr},i}, b^{\text{tr},i})\}_{i=1}^{N_{\text{tr}}}$$

$$Z_v = \{t(x^{\text{v},i}, y_1^{\text{v},i}, y_2^{\text{v},i}, b^{\text{v},i})\}_{i=1}^{N_v}$$

Density ratio estimator $w = \omega(Z_{\text{tr}}, Z_v)$

v. Choices of transformation function

Loss (IW-DPO-L)

$$t: (x, y_1, y_2, b) \mapsto \ell_{\text{DPO}}(x, y_1, y_2, b)$$

$$\ell_{\text{DPO}}(x, y_1, y_2, b) = -\log \sigma(b \cdot (r(x, y_1) - r(x, y_2)))$$

Reward (IW-DPO-R)

$$t: (x, y_1, y_2, b) \mapsto \hat{r}(x, y_1, y_2, b)$$

$$\hat{r}(x, y_1, y_2, b) = (r(y_1), r(y_2))$$

$$r(x, y) = \beta \log \frac{\pi_\theta(y | x)}{\pi_{\text{ref}}(y | x)}$$

IV. Experimental Scenarios

We simulated three deployment distribution shift scenarios

Training data
Helpful-Harmful responses
+
Helpful-Harmless responses

Training data
Science fiction-domain prompts
+
Science-domain prompts

Test data
Helpful-Harmless responses

Test data
Science-domain responses

Helpful-Harmless LM
Shift type: d or h
Dataset: SafeRLHF [3]

Science LM
Shift type: b or f
Dataset: SHP [4]

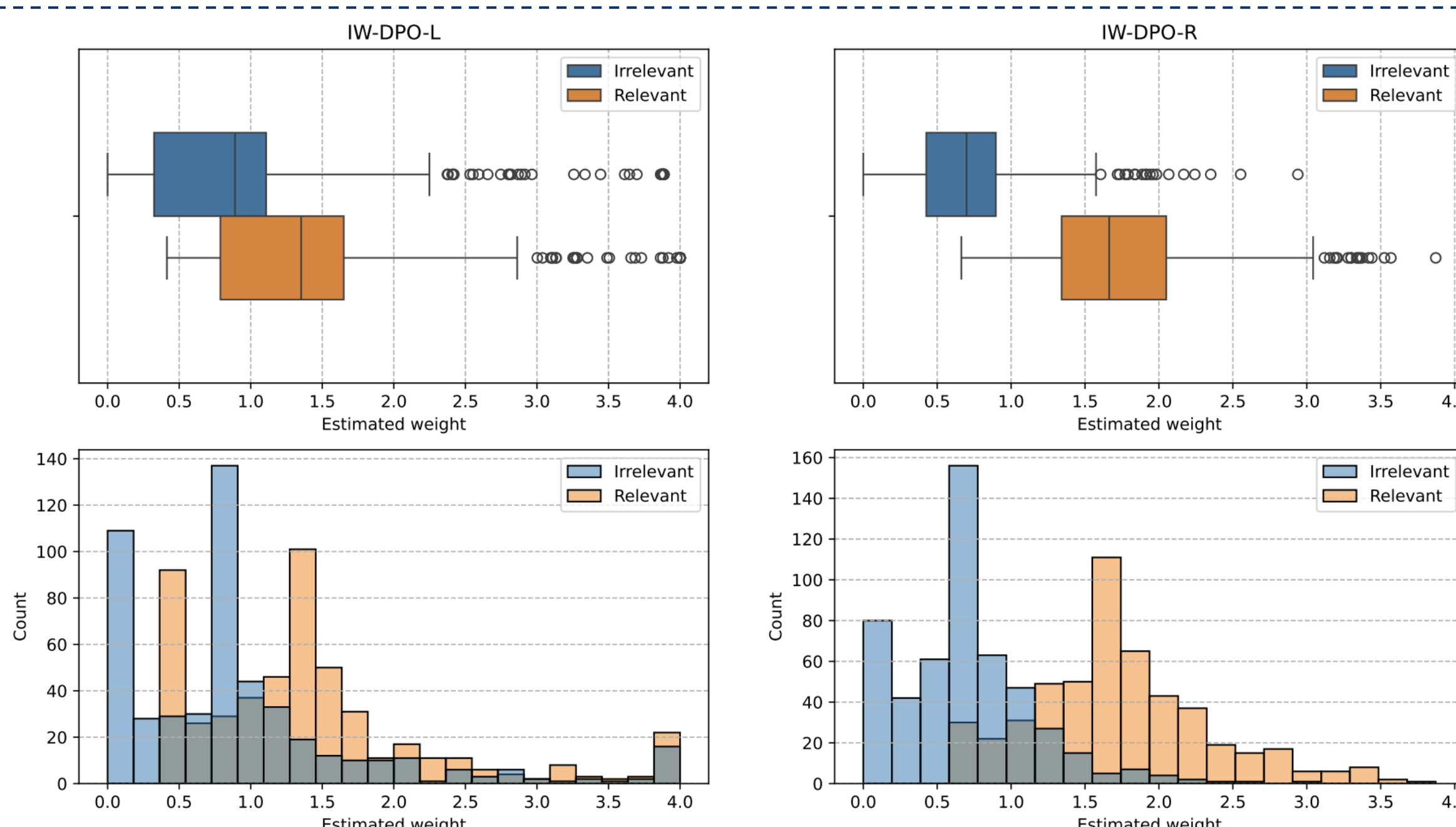
Training data
American-culture preference labels
+
Indian-culture preference labels

Test data
Indian-culture preference labels
Culture-Aware LM
Shift type: e
Dataset: CALI [5]

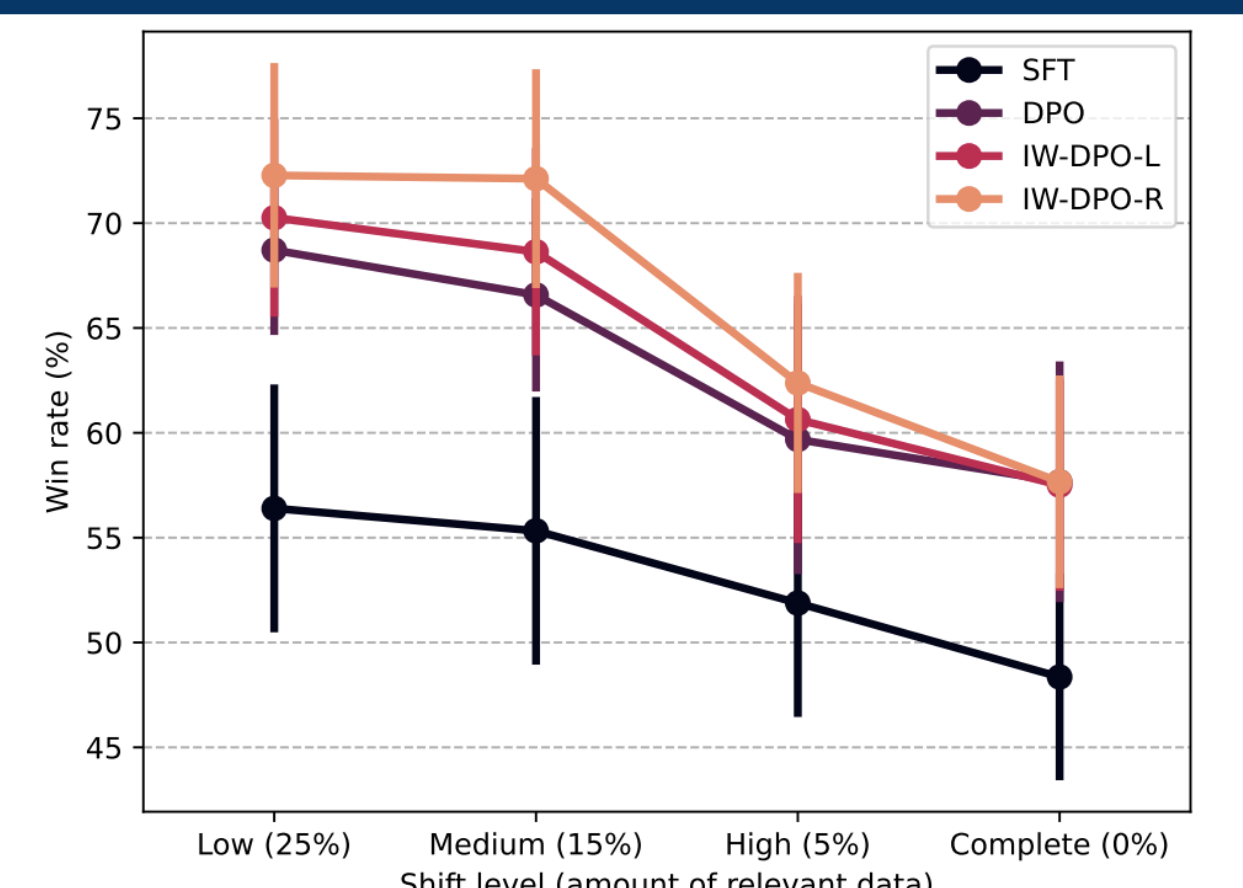
V. Results

Method	Helpful-Harmless LM	Science LM	Culture-Aware LM
SFT w/ $D_{\text{tr}} + D_v$	56.40 ± 5.12	47.06 ± 5.59	31.72 ± 3.13
DPO w/ D_v	60.48 ± 4.25	53.20 ± 5.14	32.15 ± 3.56
DPO w/ $D_{\text{tr}} + D_v$	68.71 ± 3.45	63.79 ± 3.45	35.62 ± 0.97
WPO (Zhou et al., 2024) w/ $D_{\text{tr}} + D_v$	70.26 ± 4.05	64.84 ± 5.22	36.41 ± 1.25*
IW-DPO-L	70.50 ± 3.46	65.88 ± 6.96*	36.49 ± 1.39*
IW-DPO-R	72.28 ± 4.62	70.59 ± 3.01	36.92 ± 1.77

i. **More** improvement in win rate in the **Helpful-Harmless LM** and **Science LM** scenarios, but **less** in the **Culture-Aware LM** scenario



ii. Importance weight differences become **clearer** with IW-DPO-R!



iii. Performance **degrade** when distribution shift becomes **more severe**

Scenario	Method	Density ratio estimator	Win/Match rate (%)
Helpful-Harmless LM	IW-DPO-L	KMM KLIEP RuLSIF	70.50 ± 3.46* 70.10 ± 4.39 72.28 ± 4.94
	IW-DPO-R	KMM KLIEP RuLSIF	72.28 ± 4.62* 71.88 ± 4.20* 73.19 ± 3.39
Science LM	IW-DPO-L	KMM KLIEP RuLSIF	65.88 ± 6.96* 68.10 ± 2.66 67.58 ± 3.32*
	IW-DPO-R	KMM KLIEP RuLSIF	70.59 ± 3.01 69.28 ± 4.45* 70.59 ± 4.68*
Culture-Aware LM	IW-DPO-L	KMM KLIEP RuLSIF	36.49 ± 1.39* 37.83 ± 2.68 36.45 ± 0.70*
	IW-DPO-R	KMM KLIEP RuLSIF	36.92 ± 1.77* 36.25 ± 1.36* 38.38 ± 1.46

iv. Choice of density ratio estimator is **not significant**

References

- [1] T. Fang et al. Rethinking Importance Weighting for Deep Learning under Distribution Shift. In NeurIPS, 2020.
- [2] R. Rafailov et al. Direct Preference Optimization: Your Language Model is Secretly a Reward Model. In NeurIPS, 2024.
- [3] J. Ji et al. BeaverTails: Towards Improved Safety Alignment of LLM via a Human Preference Dataset. In NeurIPS, 2023.
- [4] K. Ethayarajah et al. Understanding Dataset Difficulty with V-Usable Information. In ICML, 2022.
- [5] J. Huang et al. Culturally Aware Natural Language Inference. In EMNLP, 2023.



大学共同利用機関法人 情報・システム研究機構
統計数理研究所
The Institute of Statistical Mathematics



東京大学
THE UNIVERSITY OF TOKYO